

Gradient Descent.

$$f : \mathbb{R}^n \rightarrow \mathbb{R}$$

we need to

$$\min_{x \in \mathbb{R}^n} f(x)$$

we find the min iteratively by starting at init $x \Rightarrow x^0$

$$x_{k+1} = x_k + \alpha_k d_k$$

d_k : direction of improvement.

α_k : step size

d is direction of improvement if

$$\exists \alpha > 0 \text{ s.t. } f(x + \alpha d) < f(x)$$

So gradient descent

$f : \mathbb{R}^n \rightarrow \mathbb{R}$ differentiable then

$(x_n)_{n \in \mathbb{N}}$ sequence s.t $f(x_0) \geq f(x_1) \geq f(x_2) \dots$

so this sequence is obtainable by

$$x_{k+1} = x_k + \alpha d_k$$

where $d_k = -\nabla f(x_k)$ $\alpha_k > 0$

$$\alpha_k = \min_{\alpha_k \in \mathbb{R}_n} f(x_k - \alpha_k \nabla f(x_k))$$

Subgradient Method

if $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is not differentiable / continuous

$$x_{k+1} = x_k + \alpha_k g_k$$

g_k is any subgradient of f at x_k

α_k pre determined value. before the run of the method.

S.M will converge slowly.

maybe α_k might increase the value of f

$$\text{so } f_{k+1} = \min(f_k, f(x_k))$$

Gradient descent for quadratic form:

$f: \mathbb{R}^n \rightarrow \mathbb{R}$ $A \in \mathbb{R}^n \times \mathbb{R}^n$ symmetric
 $b \in \mathbb{R}^n$ $c \in \mathbb{R}$

$$f(x) = \frac{1}{2} x^T A x - x^T b + c$$

we know

$$\nabla f(x) = \frac{1}{2} Ax + \frac{1}{2} A^T x - b$$

$$\nabla f(x) = Ax - b$$

$$H_f(x) = A \rightarrow \text{Positive definite}$$

this means x^* $Ax^* - b = 0$ is a unique minimizer of f .

so $x^* = A^{-1}b$ we can't inverse A

we can compute α_k

$$x_{k+1} = x_k - \alpha_k \nabla f(x_k)$$

$$\alpha_k = \frac{d_k^T d_k}{d_k^T A d_k}$$

Example

$$A = \begin{pmatrix} 7.75 & 3.9 \\ 3.9 & 3.25 \end{pmatrix} \quad b = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \quad c = 0$$

$$f(x) = \frac{1}{2} x^T A x \quad \nabla f(x) = Ax$$

$$H_f(x) = A$$

Starting $x_0 = \begin{pmatrix} 1 \\ -0.8 \end{pmatrix}$ $\nabla f(x_0) = A \begin{pmatrix} 1 \\ -0.8 \end{pmatrix}$

$$x_1 = x_0 - \alpha_0 \nabla f(x_0)$$

$$x_1 = \begin{pmatrix} 1 \\ -0.8 \end{pmatrix} - 0.01 \begin{pmatrix} 11.63 \\ 1.3 \end{pmatrix} \Rightarrow \alpha_0 = 0.01$$

$$\begin{pmatrix} 1 \\ -0.8 \end{pmatrix} - \begin{pmatrix} 0.0463 \\ 0.013 \end{pmatrix} = \begin{pmatrix} 0.9537 \\ -0.813 \end{pmatrix} = x_1$$

$$\text{or } \alpha_0 = \frac{\nabla f(x_0)^T \nabla f(x_0)}{\nabla f(x_0)^T A \nabla f(x_0)} = \frac{23.1269}{228.716675}$$

$$\boxed{\alpha_0 = 0.10112}$$

better jump
 $\begin{pmatrix} 0.531 \\ -0.8013 \end{pmatrix} = x_1$

the gradient is always orthogonal to iso lines and show zigzagging characteristic.

Conjugated Gradient.

A-Orthogonal

$x, y \in \mathbb{R}^n$ are A-orthogonal or conjugated
energy norm $\langle x, y \rangle_A = 0$

$$\|x\|_A = \sqrt{\langle x, x \rangle_A} = \sqrt{x^T A x}$$

So in CG the idea is to choose the
 d_k A orthogonal to be A-orthogonal to
each other.

the idea the same

$$x_{k+1} = x_k + \alpha_k d_k$$

$\forall i, j \mid i \neq j \quad \langle d_i, d_j \rangle_A = 0 \quad$ A-orthogonal or
conjugated.

$$\alpha_k = \min_{\alpha_k > 0} f(x_k + \alpha_k d_k)$$

Residual

$$r_k = b - Ax_k$$

error step

error

$$e_k = x^* - x_k$$

\leftarrow the min

$$\text{So } r_k = b - Ax_k = Ax^* - Ax_k = A(x^* - x_k)$$

$$r_k = Ae_k$$

$$e_{K+1} = x_{K+1} - x^*$$

$$e_{K+1} = x_K + \alpha_K d_K - x^*$$

$$e_{K+1} = e_K + \alpha_K d_K$$

like Gradient descent

$$\alpha_K = \frac{d_K^T r_K}{d_K^T A d_K}$$

CG

can solve the minimization problem
in n steps giving the

$A \in \mathbb{R}^{n \times n}$ Positive definite & symmetric.
and

$$f(x) = \frac{1}{2} x^T A x - x^T b + c$$

in this quadratic form.

this means $e_n = 0$ $r_n = 0$

idea) basically A change the space

to (contour line to circles) so the gradient will go to max length

energy scalar product to change space.
to circles.

Variants of CG:

CG is restricted to A s.p.d and f as quadratic function.

Problem if A is not [S.P.D]

- apply CG on normal equation

if $A \in \mathbb{R}^{m \times n}$ $m > n$

we have $\min_{x \in \mathbb{R}^n} \frac{1}{2} \|Ax - b\|_2^2$

normal equation $A^T A x = A^T b$

$$\nabla f(x) = A^T (Ax - b)$$

$$H_f(x) = A^T A \quad \begin{matrix} \text{Positive} \\ \text{definite} \end{matrix}$$

if is not linear Function.

Solution : none linear CG / no convergence guarantees ----

we have to accelerate convergent of CG,
we have huge matrix. so like $\text{Im } \boxed{n}$

Solution : Preconditioning

the first steps get the algo a lot further towards the solution.

$$\text{we solve } M^{-1}A\bar{x} = M^{-1}\bar{b}$$

$$\text{ideally } M = A \Rightarrow x = A^{-1}b$$

$$M = \text{diag}(A) \text{ so } M^{-1} = \begin{pmatrix} \frac{1}{a_{11}} & & \\ & \frac{1}{a_{22}} & \\ & & \frac{1}{a_{33}} \end{pmatrix}$$

Tikhonov Regularization:

Let $Ax = b$ $A \in \mathbb{R}^{m \times n}$ $b \in \mathbb{R}^m$
 $\zeta \in \mathbb{R}^n$ error

$$Ax + \zeta = b \quad \Rightarrow \quad \text{usually there is error}$$

ill conditioning: $m=n$ A is invertible

$$\begin{aligned} A^{-1}(Ax + \zeta) &= A^{-1}b \\ x + A^{-1}\zeta &= A^{-1}b \\ x &= A^{-1}b - \underbrace{A^{-1}\zeta}_{\text{(err)}} \quad (\text{err}) \end{aligned}$$

estimating the error size $\|A^{-1}\zeta\|_2 \leq \|A^{-1}\|_2 \|\zeta\|_2$
 this
 is
 small

matrix norm

$$\|A\|_2 = \sup_{\substack{x \neq 0 \\ \text{max}}} \frac{\|Ax\|_2}{\|x\|_2}$$

or $\|A\|_2 = \sigma_{\max}(A) \quad \text{max singular value.}$

$$\|Ax\| \leq \|A\| \|x\|$$

Condition of matrix: $A \in \mathbb{R}^{m \times n}$

$$\text{cond}(A) = \frac{\sigma_1(A)}{\sigma_n(A)} \leftarrow \begin{array}{l} \text{singular} \\ \text{value.} \end{array}$$

$$\sigma_1 > \sigma_2 > \dots > \sigma_n$$

$$\text{cd} \quad \sigma_1 = 1 \quad \sigma_n \lll 1$$

$$\text{so } \text{cond}(A) = \sigma_{\max}(A)$$

$$\text{so for } \text{cond}(A^{-1}) = \frac{1}{\sigma_n} \underset{\leftarrow \text{ of } A}{>>>}$$

large error so the issue in A

So $\|A^{-1}\|$ is huge
in Tikhonov to reduce the effect

to noise (explosion of errors).

So the biggest σ in $A = U \Sigma V^T$

$$\hookrightarrow A^{-1} = (U \Sigma V^T)^{-1} = V \Sigma^{-1} U^T$$

$$\Sigma = \begin{pmatrix} \sigma_1 & & \\ & \sigma_2 & \\ & & \ddots & \sigma_n \end{pmatrix}$$

$$\Sigma^{-1} = \begin{pmatrix} \frac{1}{\sigma_1} & & 0 \\ & \frac{1}{\sigma_2} & \dots & 0 \\ 0 & & & \frac{1}{\sigma_n} \end{pmatrix} \text{ this is not } U \Sigma V^T \text{ or } A^{-1}$$

$$\text{but they are the same in } A^{-1} = V \underline{\Sigma^{-1}} U^T$$

not SVD

Cope with huge $\text{cond}(A^{-1})$

use truncated SVD remove all small singular values.

$$A = U \Sigma V^T \quad \Sigma = \text{diag}(\sigma_1, \dots, \sigma_n)$$

$$\text{approx. } \hat{x} = G_K(b) = V \Sigma_K^+ U^T b$$

$$\text{where } \Sigma_K^+ = \text{diag}\left(\frac{1}{\sigma_1}, \frac{1}{\sigma_2}, \dots, \frac{1}{\sigma_K}, 0, 0, 0\right)$$

$$\text{we have } G_K(b) = V \Sigma_K^+ U^T (Ax + \varepsilon)$$

$$= V \Sigma_K^+ U^T U \Sigma V^T x + V \Sigma_K^+ U^T \varepsilon$$

$$= V \Sigma_K^+ \underbrace{\Sigma V^T x}_{\text{error}} + V \Sigma_K^+ U^T \varepsilon$$

$$\leq \|V \Sigma_K^+ U^T\| + \|\varepsilon\| \quad \leftarrow \text{small}$$

V and U^T are orthogonal and length preserving so we drop them

$$\| \Sigma_K^+ \|_2 = \sigma_{\max} = \frac{1}{\sigma_K}$$

$$\text{so the error is } = \frac{1}{\sigma_K} \|\varepsilon\|$$

$$\text{since } \frac{1}{\sigma_1} < \frac{1}{\sigma_2} < \frac{1}{\sigma_K} \dots$$

more we remove the influence of ε but also we reduce the solution correctness.

use Tikhonov regularization to reduce the error influence.

Tikhonov Regularized solution.

$$T(b) = \min_{\alpha} \|A\bar{z} - b\|^2 + \alpha \|z\|_2^2$$

for $\alpha > 0$ regularized parameter.

$$\left. \begin{array}{l} A = U \Sigma V^T \\ T_\alpha(b) = V D_\alpha^+ V^T b \end{array} \right\} \begin{array}{l} D_\alpha^+ = \text{diag}\left(\frac{\sigma_1}{\sigma_1^2 + \alpha}, \dots, \frac{\sigma_n}{\sigma_n^2 + \alpha}\right) \\ \text{make } \alpha \text{ bigger then} \\ \text{make the quotient smaller} \\ \text{so it doesn't explode the error.} \end{array}$$

Stack Form.

$$\begin{bmatrix} A \\ \sqrt{\alpha} I \end{bmatrix} x = \begin{bmatrix} b \\ 0 \end{bmatrix}$$

$$\tilde{A} \tilde{x} = \tilde{b}$$

$$\tilde{A}^T \tilde{A} \tilde{x} = \tilde{A}^T \tilde{b}$$

$$\begin{bmatrix} A^T \sqrt{\alpha} I \\ \sqrt{\alpha} I \end{bmatrix} \begin{bmatrix} A \\ \sqrt{\alpha} I \end{bmatrix} = A^T A + \alpha I$$

$$(A^T A + \alpha I)x = \tilde{A}^T \tilde{b}$$

$$\begin{aligned}\tilde{A}^T \tilde{b} &= [A^T \ \sqrt{\alpha} I] \begin{bmatrix} b \\ 0 \end{bmatrix} \\ &= A^T b\end{aligned}$$

so

$$(A^T A + \alpha I)x = A^T b$$

$$x = (A^T A + \alpha I)^{-1} A^T b$$

if A is full rank $A^T A$ is

Symmetric positive definite we can use

CG to solve this.

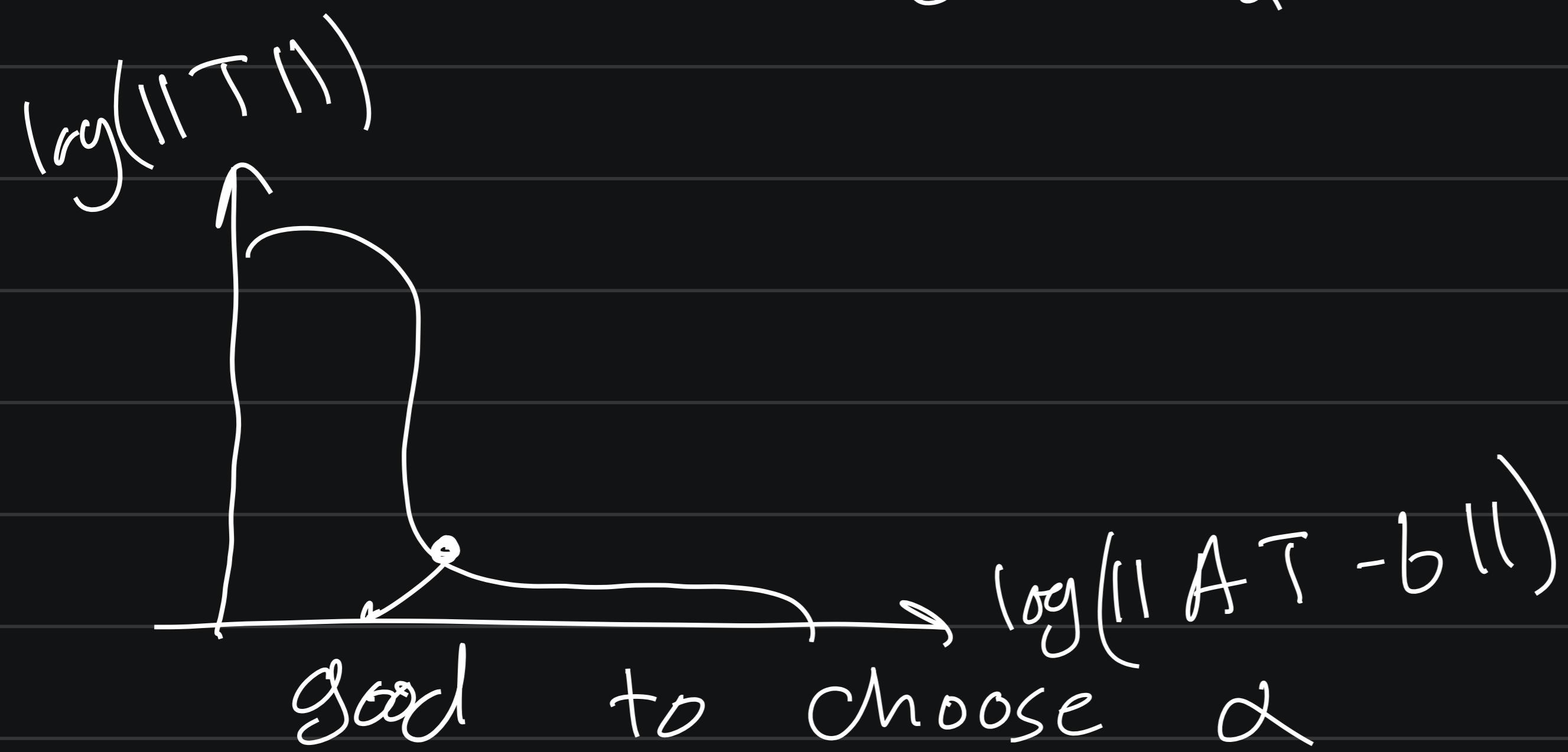
how to choose α

L curve method

we choose candidate of α

we compute the solution for each α

2D curve $(\log \|A^T \tau_\alpha(b) - b\|, \log \|\tau_\alpha(b)\|)$



Newton Method

Solve system

$$\Omega \subset \mathbb{R}^n, F: \Omega \rightarrow \mathbb{R}^n$$

of non-linear Function

Solve

$$F(x) = 0$$

setting F to 0

$$\text{start } x_0 \in \Omega \Rightarrow x_{k+1} = x_k - J_F^{-1}(x_k) F(x_k)$$

Guess Newton Method (minimizing)

$$g: \mathbb{R}^n \rightarrow \mathbb{R}^m \text{ solve } \min_{x \in \mathbb{R}^n} \frac{1}{2} \|g(x)\|_2^2$$

$$\text{start } x_0 \in \Omega \Rightarrow x_{k+1} = x_k - (J_g(x_k) J_g(x_k))^{-1} J_g(x_k)^T g(x_k).$$

Fixed Point

Let $N \subset M \neq \emptyset$ $T: M \rightarrow N$

a solution x^* is called fixed point if

$$T(x^*) = x^* \text{ doesn't change.}$$

example

- $T(x) = x + 1 \rightarrow$ no fixed points.

f.p.s does not converge

- $T_2(x) = ax \quad a \in R \neq 1$

if $a \neq 1$ only $x=0 \quad T(0) = a0 = 0$
so $x=0$ fixed point.

f.p.t converges to 0 if $|a| < 1$ and
doesn't not converge $|a| > 1$

- $T_3(x) = x \quad$ every x is a fixed point.

convergence at any x

Fixed Point iteration

(M, d) $T: M \rightarrow M$ iteration

$$x_0 \in M \quad x_{k+1} = T(x_k)$$

will it lead to a fixed point

we will stop when $x_{k+1} = T(x_k) = x_k$

Lipschitz-continuous / contraction

(M, d) $A \subset M$ $T: A \rightarrow M$ if $\exists L > 0$

s.t $d(T(x), T(y)) \leq L \cdot d(x, y)$
 $\forall x, y \in A$

if $L \in (0, 1)$ then we call it contraction.

Banach fixed point theorem.

Let (X, d) $M \subset X$ $M \neq \emptyset$ $T: M \rightarrow M$

- ① Complete space imp
- ② M closed
- ③ T is proper contraction

then $\exists x^*$ s.t. (exactly one) s.t
 $x^* \in M$ and fixed point iteration
converges to it.

$$\text{then } d(x_k, x^*) \leq \frac{L}{1-L} d(x_{k-1}, x_k)$$

Brouwer

Let $M \subset \mathbb{R}^n$

- ① Convex
 - ② Compact
 - ③ $\neq \emptyset$
- $T: M \rightarrow M$ continuous
- then there is at least one fixed point x^*

Fixed Point iteration to solve non-linear equation

$$F(x^*) = 0 \Leftrightarrow T(x^*) = x^*$$

$$\text{So } T(x_1) = x - J_F(x)^{-1} F(x).$$

important

fixed point iteration : everything depends on T

typical : $T(x) = x - F(x)$

so if $T(x^*) = x^*$

① then $T(x^*) = x^* - F(x^*)$

then $F(x^*) = 0 \Rightarrow$

② $T(x) = x - J_F^{-1}(x) F(x)$ newton method

F.P.I is the most generic type

of optimization.

General Optimization Problem

$$\min_x D(x, y) + \lambda R(x)$$

we can add multiple Regularization + $\sum_{i=1}^n \lambda_i R_i(x)$

important question

if we see some optimization we need to choose which tool we can use for it to optimize it

for example if we have $\|x\|_1$, \leftarrow not differentiable

then GD, Does not work
or newton doesn't also
may be Subgradient method works.